



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### Model updating after interventions paradoxically introduces bias

**Citation for published version:**

Liley, J, Emerson, SR, Mateen, BA, Vallejos, CA, Aslett, LJM & Vollmer, SJ 2021, 'Model updating after interventions paradoxically introduces bias', Paper presented at 24th International Conference on Artificial Intelligence and Statistics, 13/04/21 - 15/04/21.

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



---

# Model updating after interventions paradoxically introduces bias

---

Anonymous Author  
Anonymous Institution

## Abstract

Machine learning is increasingly being used to generate prediction models for use in a number of real-world settings, from credit risk assessment to clinical decision support. Recent discussions have highlighted potential problems in the updating of a predictive score for a binary outcome when an existing predictive score forms part of the standard workflow, driving interventions. In this setting, the existing score induces an additional causative pathway which leads to miscalibration when the original score is replaced. We propose a general causal framework to describe and address this problem, and demonstrate an equivalent formulation as a partially observed Markov decision process. We use this model to demonstrate the impact of such ‘naive updating’ when performed repeatedly. Namely, we show that successive predictive scores may converge to a point where they predict their own effect, or may eventually oscillate between two values, and we argue that neither outcome is desirable. Furthermore, we demonstrate that even if model-fitting procedures improve, actual performance may worsen. We complement these findings with a discussion of several potential routes to overcome these problems.

## 1 Introduction

A common machine learning task concerns the prediction of an outcome  $Y$  given a known set of predictors  $X$  [Friedman et al., 2001]. Usually, the intent is to anticipate the value of  $Y$  in situations in which only  $X$  is known. Often, the ultimate goal is to avoid or en-

courage certain values of  $Y$ , with interventions guided by the predictions provided by the algorithm.

We focus on the standard setting, often seen in health-care, where  $X$  is first observed and used to make predictions about  $Y$ , then interventions occur before outcomes are observed. This setting can lead to prediction scores being ‘victims of their own success’ [Lenert et al., 2019, Sperrin et al., 2019]. Interventions driven by the score can change the distribution of the data and outcomes, leading to a decay in observed performance, particularly if the intervention is successful. Analysis of this effect requires consideration of the causal processes governing  $X$ ,  $Y$ , and the potential interventions driven by the score [Sperrin et al., 2019]. Predictive scores are often implemented by direct dissemination to agents that are capable of modifying these causal processes [Rahimian et al., 2018, Hyland et al., 2020], which leads to vulnerability to this problem. This problem also exist if predictions influence discrete actions, initial progress for this has been made using bandits [Shi et al., 2020].

This problem is particularly critical in settings where existing predictive scores are to be replaced by an updated version. In many real-world contexts, the underlying phenomena represented by the predictive model will change over time [Wallace et al., 2014]; statistical procedures for prediction may also improve (particularly for complex tasks); and researchers may wish to include further predictors or increase the scope of predictive scores. In general, we may expect that most predictive algorithms will need to be updated or replaced over time. Up-to-date models should generally be trained on the most recent available data which, as described above, will be contaminated by interventions based on existing scores. Should a new predictive model be fitted to new observations of  $X$  and  $Y$ , it will consequently also model the impact of the existing score. Removal of the existing score will introduce bias into predictions made by the new score, as will insertion of the new score in place of the old. We term such an operation a ‘naive model replacement’.

tively studied. We use this framework to draw attention to the hazards of naive model replacement, especially when it occurs repeatedly. We introduce these hazards in the context of a generalised ultimate aim of the model, formulated as a constrained optimisation problem in which the occurrence of undesirable values of  $Y$  is to be minimised with limited intervention.

A simple parable of this phenomenon concerns yearly influenza vaccinations. In a vaccination-naïve population, risk assessments for influenza motivate widespread vaccination. However, in a later ‘epoch’, the risk may appear much lower, and could naively suggest vaccination is no longer required introducing risks to public health<sup>1</sup>. More generally, updated risk scores for clinical outcomes may be biased due to the interventions motivated by the scores themselves. As a second example, consider risk scores used to predict future emergency hospital admissions  $Y$ , on the basis of covariates  $X$  [Rahimian et al., 2018]. Suppose that prescription of some drug  $D \in X$  confers increased risk, and this is established by the risk score. Should such risk scores be distributed at time  $t = 0$  to agents able to modify these factors (e.g., doctors), they may intervene by taking patients off  $D$  thereby reducing emergency admission risk  $\mathbb{E}[Y]$  at a time  $t = 1$ . If a new score is naively fitted to  $X$  at  $t = 0$  and  $Y$  at  $t = 1$ , it would underestimate the danger of  $D$ .

Section 2 describes the problem in terms of causal effects. We develop this into a full model specification in Section 2.2, along with a description of the constrained optimisation problem the model/intervention pair aims to solve in 2.3. In Section 3, we analyse the short and long-term effects of repeated naive replacement and show that they are generally undesirable. In Section 4, we discuss three classes of solutions: more complex modelling, routine maintenance of a ‘hold-out’ set, and controlled interventions. In Section 5 we describe a reformulation of the model as control theory problem. Finally, in Section 6, we discuss limitations and implications of our approach. Our supplementary material contains relevant examples and proofs, an exposition of the problem in a real-world example, and a list of open problems in this setting.

## 2 Model

### 2.1 Overview

Assume that we are attempting to predict an outcome  $Y$  given a known set of covariates  $X$ . For simplicity, we assume  $Y$  is a binary (e.g. admission versus non admission to an Intensive Care Unit) and model it as

<sup>1</sup>See for example <https://www.who.int/news-room/spotlight/ten-threats-to-global-health-in-2019>

a Bernoulli random variable. If  $Y = 1$  is considered to be a negative outcome, often the eventual aim is to reduce  $\mathbb{P}(Y = 1|X) = \mathbb{E}[Y|X]$ ; we will discuss this in Section 2.2 once we have defined terms formally. For the moment, we assume the causal structure shown in Figure 1. We denote by  $\rho_0(X)$  an initial predictive model for  $\mathbb{E}[Y|X]$ , fitted to observations of  $(X, Y)$  generated under the causal structure in Figure 1A. During deployment, we compute  $\rho_0(X)$  for all members of a population and disseminate it to *agents who can intervene* on  $X$  (e.g. doctors) based on those predictions, aiming to prevent  $Y = 1$ . Replacing or updating  $\rho_0$ , will typically involve fitting a new predictive model  $\rho_1(X)$  to new observations of  $(X, Y)$ . It is clear that while  $\rho_0(X)$  is an estimator of  $\mathbb{E}[Y|X]$ , the new predictive function  $\rho_1(X)$  is instead an estimator of

$$\mathbb{E}[Y|X, \text{do}[\rho_0(X)]] \quad (1)$$

where  $\text{do}[\rho_0(X)]$  indicates the action ‘compute and disseminate  $\rho_0(X)$ ’. Although  $\rho_0(X)$  is determined by  $X$ , the computation  $\text{do}[\rho_0(X)]$  makes  $\rho_0$  actionable. This opens a second causal pathway from  $X$  to  $Y$ , affecting the setting in which  $\rho_1$  is fitted (Figure 1B). If the initial score  $\rho_0(X)$  is universally disseminated, the distribution of  $Y$  given  $X$  (without the  $\text{do}[\rho_0(X)]$ ) now becomes a counterfactual which we cannot observe.

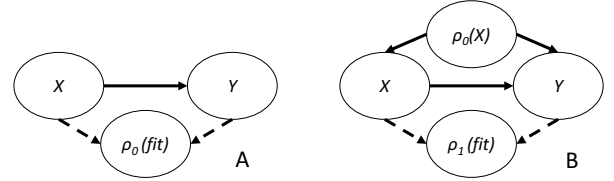


Figure 1: Causal structure under which  $\rho_0$  (panel A) and  $\rho_1$  (panel B) are fitted. Dashed lines indicate a model-fitting process.

### 2.2 General notation and assumptions

Here, we use a causal model to illustrate potential emergent behaviour resulting from repeated naive model updating. We do not aim to cover the complexities of *all* real-world applications, yet our simplified setup is sufficient to demonstrate the dangers arising in this context.

As  $\rho_0$  is deployed and drives interventions, covariate values  $X$  may change, as may the dependence of  $Y$  on  $X$ . Here, we partition  $X$  into three sets:

- $X^s$ : Fixed or ‘set’ covariates;  $\dim(X^s) = p^s$ ,
- $X^a$ : Actionable covariates;  $\dim(X^a) = p^a$ ,
- $X^\ell$ : Latent covariates;  $\dim(X^\ell) = p^\ell$ . (2)

Although  $X^\ell$  may influence the causal mechanism between  $X$  and  $Y$  and may be intervened on, we assume it is unobserved. Hence, only  $X^s$  and  $X^a$  are known when evaluating a risk score, and  $X^s$  cannot be intervened on (e.g. ‘Age’). We also define two sets of time indicators  $t, e$  (time, epoch):

$$t \in \{0, 1\} : \begin{cases} t = 0: \text{predictive score is computed} \\ t = 1: Y \text{ observed, after possible} \\ \quad \text{intervention} \end{cases}$$

$$e \in \mathbb{N} : \begin{cases} e = 0: \text{no predictive score is used} \\ e > 0: \text{model from epoch } e - 1 \text{ is used.} \end{cases}$$

We assume that values of  $X$  depend on  $t$  and  $e$  using the notation  $X_e(t) = (X_e^s(t), X_e^a(t), X_e^\ell(t)) \in \Omega^s \times \Omega^a \times \Omega^\ell = \Omega$ . As  $Y$  is only observed at  $t = 1$ ,  $Y$  at epoch  $e$  is denoted as  $Y_e$ . At each epoch, we assume that values of  $X_e(t)$  across individuals in the population are *iid* with probability measure  $\mu_e$ . We introduce the following functions

$$\begin{aligned} f_e(x^s, x^a, x^\ell) &= \mathbb{E}[Y_e | X_e(1) = (x^s, x^a, x^\ell)] \\ &= \text{Causal mechanism determining} \\ &\quad \text{probability of } Y_e = 1 \text{ given } X_e(1) \\ g_e^a(\rho, x^a) &\in \{g : [0, 1] \times \Omega^a \rightarrow \Omega^a\} \\ &= \text{Intervention process on } X^a \text{ in} \\ &\quad \text{response to a predictive score } \rho \\ &\quad \text{updating } X_e^a(0) \rightarrow X_e^a(1) \\ g_e^\ell(\rho, x^\ell) &\in \{g : [0, 1] \times \Omega^\ell \rightarrow \Omega^\ell\} \\ &= \text{Intervention process on } X^\ell \text{ in} \\ &\quad \text{response to a predictive score } \rho \\ &\quad \text{updating } X_e^\ell(0) \rightarrow X_e^\ell(1) \\ \rho_e(x^s, x^a) &\in \{\rho_e : \Omega^s \times \Omega^a \rightarrow [0, 1]\} \\ &= \text{Predictive score trained at epoch} \\ &\quad e, \text{ evaluated at observed covariates.} \end{aligned}$$

Our main model is based on the following assumptions

1.  $\forall e \ X_e^s(0) = X_e^s(1)$ : ‘set’ covariates do not change from  $t = 0$  to  $t = 1$
2.  $X_0^a(0) = X_0^a(1)$ ,  $X_0^\ell(0) = X_0^\ell(1)$ : ‘actionable’ and ‘latent’ covariates do not change at epoch 0
3.  $X_e^\ell(t)$  is unobserved, but may be modified from  $t = 0$  to  $t = 1$  in response to  $\rho_{e-1}$
4. Values of  $X_e(0)$  are independent across epochs, i.e. we do not track the same subjects over time.
5. At epoch  $e$ , the predictive score uses only  $X_e^a(0)$ ,  $X_e^s(0)$  and  $Y_e$  as training data; previous epochs are ignored and  $X_e^a(1)$ ,  $X_e^s(1)$  are not observed.

6.  $\forall e \ \mathbb{E}[Y_e | X_e] = \mathbb{E}[Y_e | X_e(1)]$ :  $Y_e$  depends only on  $X_e(1)$ ; that is, after any potential interventions.

Besides these core assumptions, for the applications in this work, we variably assume some of the following

7.  $f_e$ ,  $g_e^a$ ,  $g_e^\ell$  and  $\mu_e$  remain fixed across epochs<sup>2</sup>, so values  $\{X_e^s\}$  are *iid*, as are  $\{X_e^a\}$  and  $\{X_e^\ell\}$  (within an epoch they may be correlated). Where we make this assumption, we will omit the epoch subscript for clarity. We also use the shorthand  $X^\ell \equiv X_e^\ell(0) | (X_e^s(0), X_e^a(0)) = (x^s, x^a)$
8. We allow  $\rho_e$  to be an arbitrary function, but generally presume it is an estimator of

$$\begin{aligned} \rho_e(x^s, x^a) &\approx \mathbb{E}[Y_e | X_e^s(0) = x^s, X_e^a(0) = x^a] \\ &= \mathbb{E}_{X^\ell} [f_e(x^s, g_e^a(\rho_{e-1}, x^a), g_e^\ell(\rho_{e-1}, X^\ell))] \\ &\triangleq \tilde{f}_e(x^s, x^a) \end{aligned} \quad (3)$$

noting that  $\tilde{f}_e$  depends on  $e$  even if  $f_e$  does not.

9. The function  $f_e$  is  $C^1$  in all arguments, and covariates are coded such that increases in covariate values increase risk
10.  $g_e^\ell$ ,  $g_e^a$  are  $C^1$  in all arguments, and a higher value of  $\rho$  means a larger intervention is made (we assume  $g_e^\ell$  and  $g_e^a$  to be deterministic, but random valued functions may more accurately capture the uncertainty linked to real-world interventions).

This extended causal model is shown in Figure 2. To aid interpretation, a real-world example is described using this notation in Supplementary Section 1.

### 2.3 Aim of predictive score

The aim of the predictive score is generally to estimate  $\mathbb{E}[Y_e | X_e(0)]$  accurately, presuming that we take  $X_e(0)$  to be identically distributed over the population concerned. However, if action is to be taken on the score, we may presume the ultimate goal is to minimise  $\mathbb{E}[Y_e]$ , i.e. minimising

$$\begin{aligned} \mathbb{E}[Y_e] &= \mathbb{E}_{X_e(0)} [Y_e | X_e(1)] \\ &= \mathbb{E}_{X_e(0)} [f_e(X^s, g_e^a(\rho, X_e^a(0)), g_e^\ell(\rho, X_e^\ell(0)))] \end{aligned} \quad (4)$$

However, we presume that we cannot afford to maximally intervene in all cases. Suppose the cost of lowering  $X^a$  and  $X^\ell$  by  $x$  is  $c^a(X^a, x)$  and  $c^\ell(X^\ell, x)$ , respectively. The total intervention must then satisfy

<sup>2</sup>In practice, we may assume  $f_e$  changes slightly between epochs, but that this change is negligible.

$$\mathbb{E}_{X_e(0)} \left[ c^a \left( X_e^a(0), X_e^a(0) - g_e^a(\rho, X_e^a(0)) \right) + c^\ell \left( X_e^\ell(0), X_e^\ell(0) - g_e^\ell(\rho, X_e^\ell(0)) \right) \right] \leq C \quad (5)$$

for a known constant  $C$ , representing maximum cost. Thus we want to minimise (4) subject to (5). We have allowed  $f_e$ ,  $\mu_e$ ,  $g_e^a$ ,  $g_e^\ell$  and  $\rho_e$  to vary across epochs. Of these, we can consider  $f_e$  and  $\mu_e$  to vary as a consequence of underlying processes, and  $g_e^a$ ,  $g_e^\ell$  and  $\rho_e$  to be (somewhat) under our control. Depending on the problem, we may either consider  $g_e^a$  and  $g_e^\ell$  as fixed, and choose an optimal function  $\rho_e$ ; or consider  $\rho_e$  as fixed, and choose optimal functions  $g_e^a$ ,  $g_e^\ell$ . If both are optimised, this corresponds to a general problem of resource allocation; see Supplementary Section 3.1

### 3 Naive model updating

We consider a ‘naive’ process in which a new score  $\rho_e$  is fitted in each epoch, and then used as a drop-in replacement of an existing score  $\rho_{e-1}$ . We show that this procedure does not generally solve the constrained optimisation problem in Section 2.3, can lead to ‘worse’ performance of ‘better’ models, and may lead to wide oscillation of predictions for fixed inputs across epochs.

#### 3.1 Worse performance of better models

Here, we show that naive updating can lead to a loss in observed performance — even when the procedure to infer  $\rho_e$  is more accurate. We adopt assumptions 1–10, taking the approximation in equation (3) to be imperfect. Although most model elements are conserved across epochs (assumption 7), we presume that the procedure used to infer  $\rho_e$  changes, leading to better estimators of the function  $\tilde{f}_e$ .

At epoch  $e$ , the training data is denoted by  $(X_e^*, Y_e^*)$  and consists of  $n$  samples of  $(X_e(0), Y_e)$ , with the latent covariate information removed. In the absence of interventions, we assert that model performance will improve over epochs. Since performance under non-intervention is equivalent to performance at epoch 0, this can be stated as:

$$\mathbb{E}_{(X_0^*, Y_0^*)} \left[ m_{\tilde{f}_0}(\rho_e | X_0^*, Y_0^*) \right] > \mathbb{E}_{(X_0^*, Y_0^*)} \left[ m_{\tilde{f}_0}(\rho_{e+1} | X_0^*, Y_0^*) \right], \quad (6)$$

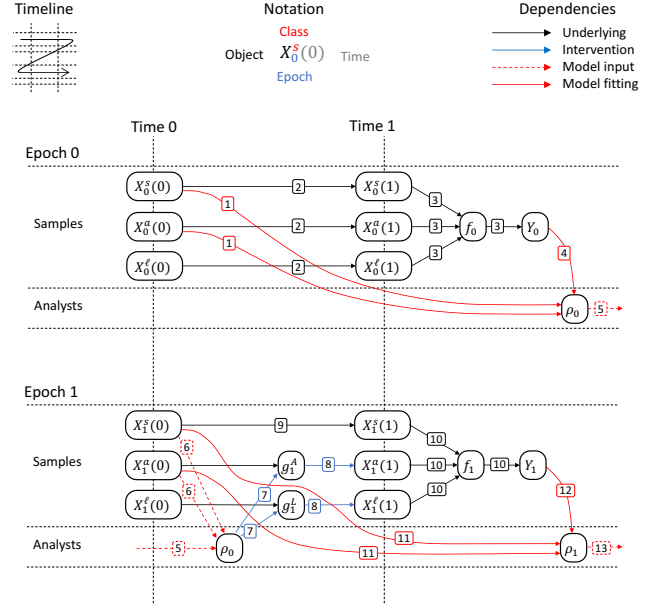


Figure 2: This figure shows a causal diagram. An ‘epoch’ is a new model fitting cycle. Covariates for a sample at the start of an epoch are modelled by  $X_e^*(0)$ . We presume  $\{X_e^s(0), e \geq 0\}$  are independent (as are  $X_e^a(0)$  and  $X_e^\ell(0)$ ). We start with a sample at  $t = 0, e = 0$ . The values  $X_0^s(0)$ ,  $X_0^a(0)$  are observed and sent to analysts (arrow 1). No predictive score is present and no interventions are made based on it, so values remain the same to  $t = 1$  (arrows 2).  $\mathbb{E}[Y_0]$  depends only on covariates at  $t = 1$ , through  $f_0$  (arrows 3).  $Y_0$  is observed and sent to analysts (arrow 4) who decide a function  $\rho_0$ , which is retained into epoch 1 (arrow 5). We start epoch 1 with a new independent sample. At  $t = 0$ , we observe  $X_0^s(0)$ ,  $X_0^a(0)$  and send them to analysts (arrow 6) who compute  $\rho_0$  ( $X_0^s(0), X_0^a(0)$ ) which is used to inform interventions  $g_1^a$ ,  $g_1^\ell$  (arrow 7) to change values  $X_e^a(0), X_e^\ell(0)$  to  $X_e^a(1), X_e^\ell(1)$  respectively (arrows 8).  $X_e^s(0)$  is not interventionable and becomes  $X_e^s(1)$  (arrow 9).  $\mathbb{E}[Y_1]$  is determined by covariates at  $t = 1$  (arrows 10). Analysts use the values of  $X_1^s(0)$ ,  $X_1^a(0)$  (arrows 11), and  $Y_1$  (arrow 12) to decide a  $\rho_1$ , which is retained (arrow 13) for epoch 2. Subsequent epochs proceed similarly to epoch 1.

where  $m_{\tilde{f}}(\rho | X, Y)$  denotes a metric for closeness of  $\rho$  to  $\tilde{f}$ , given observed data  $(X, Y)$ <sup>3</sup>. However, if interventions are in place, the improvement in equation (6), does not imply that the actual performance improves across epochs, that is:

<sup>3</sup>In practice,  $m_{\tilde{f}_e}$  is unknown but (assuming latent covariates have a small influence on  $f$ ) estimates of  $m_{\tilde{f}_0}$  can be calculated through a holdout test data set.

$$\mathbb{E}_{(X_e^*, Y_e^*)} \left[ m_{\tilde{f}_e}(\rho_e | X_e^*, Y_e^*) \right] > \mathbb{E}_{(X_{e+1}^*, Y_{e+1}^*)} \left[ m_{\tilde{f}_{e+1}}(\rho_{e+1} | X_{e+1}^*, Y_{e+1}^*) \right]. \quad (7)$$

This is proven this by counterexample, see Supplementary Section 3.2. A critical consequence of this artefact is that stakeholders may decide not to update an existing score, even if an apparently better one is available.<sup>4</sup>

### 3.2 Dynamics of repeated naive updating

Here, we analyse the dynamics of repeated naive model updating. For this purpose, we make assumptions 1-10 and assume that  $\rho_e$  is an oracle: the ‘ $\approx$ ’ in equation (3) is replaced by an ‘=’.

At epoch 0, there are no interventions, hence the risk of observing  $Y = 1$  is  $\mathbb{E}[Y_0 | X_0(0) = (x^s, x^a, x^\ell)] = f(x^s, x^a, x^\ell)$ . The score  $\rho_0$  is therefore defined as

$$\rho_0(x^s, x^a) = \mathbb{E}_{X^\ell} [f(x^s, x^a, X^\ell)], \quad (8)$$

where  $X^\ell$  is denoted as in assumption 7. In subsequent epochs,  $\rho_e$  is used to modify  $x^a$  and  $x^\ell$  via  $g^a$  and  $g^\ell$ , leading to the following recursive relation:

$$\begin{aligned} \rho_0(x^s, x^a) &= \mathbb{E}_{X^\ell} [f(x^s, x^a, X^\ell)] \\ \rho_e(x^s, x^a) &= \mathbb{E}_{X^\ell} [f(x^s, g^a(\rho_{e-1}(x^s, x^a), x^a), \\ &\quad g^\ell(\rho_{e-1}(x^s, x^a), X^\ell))] \\ &\triangleq h(\rho_{e-1}(x^s, x^a)) \end{aligned} \quad (9)$$

We briefly explore the dynamics of this recursion. Let  $z \in [0, 1]$  be arbitrary and denote by  $S$  the substitution  $(x^s, x^a, x^\ell) = (x^s, g^a(z, x^a), g^\ell(z, X^\ell))$ . Recalling definitions of  $p^s, p^a$  from (2), we set (for  $i$  across the dimensions of  $(x^a, x^\ell)$ )

$$\begin{aligned} \delta_i^{g^a} &= \frac{\partial [g^a(z, x^a)]_i}{\partial z} & \delta_i^{g^\ell} &= \frac{\partial [g^\ell(z, x^\ell)]_i}{\partial z} \\ \delta_i^{f^a} &= (\nabla f|_S)_{p^s+i} & \delta_i^{f^\ell} &= (\nabla f|_S)_{p^s+p^a+i} \end{aligned}$$

recalling assumptions 9,10 to assert that these partial derivatives exist. Assumptions 9 and 10 further imply  $\delta_i^{f^L} > 0$ ,  $\delta_i^{f^A} > 0$  and  $\delta_i^{g^A} < 0$ ,  $\delta_i^{g^L} < 0$  respectively, so

$$h'(z) = \mathbb{E}_{X^\ell} \left[ \sum_i^{p^a} \delta_i^{g^A} \delta_i^{f^A} + \sum_i^{p^\ell} \delta_i^{g^L} \delta_i^{f^L} \right] < 0 \quad (10)$$

and thus the recursion  $\rho_{e+1} = h(\rho_e)$  has exactly one fixed point. Call this  $z_0$ , so  $z_0 = h(z_0)$ . We now note

<sup>4</sup>We note that practically (if a holdout test data set was used) the conclusions on performance made by stakeholders would be based on a risk score’s closeness to  $\tilde{f}_0$  instead of  $\tilde{f}_e$ , but the results are the same, which we show in Supplementary Section 3.2.

**Theorem 1.** *If  $h'(z_0) \leq -1$  then the recursion does not converge unless  $\rho_0 = z_0$ , and will converge to oscillating between two values. If for some (possibly unbounded) interval  $R$  we have  $\rho_e \in R$  for some  $n$  and for all  $z \in R$ ,  $h(z) \in R$  and*

$$\sum_i^{p^a} (\delta_i^{g^a})^2 \leq k_1, \quad \sum_i^{p^\ell} \mathbb{E}_{X^\ell} \left[ (\delta_i^{g^\ell})^2 \right] \leq k_2 \quad (11)$$

$$\sum_i^{p^a} \mathbb{E}_{X^\ell} \left[ |\delta_i^{f^a}|^2 \right] \leq k_3, \quad \sum_i^{p^\ell} \mathbb{E}_{X^\ell} \left[ (\delta_i^{f^\ell})^2 \right] \leq k_4 \quad (12)$$

where  $\sqrt{k_1 k_3} + \sqrt{k_2 k_4} < 1$ , then

$$|\rho_e(x^s, x^a) - \rho_{e+1}(x^s, x^a)| \rightarrow 0$$

as  $n \rightarrow \infty$ .

This is proved in Supplementary Appendix 3.3

Condition (11) states that, on average, interventions make only small change to  $x^a$  and  $x^\ell$  in response to small changes in  $\rho$ . Condition (12) states that, on average, the actual risk changes little with small changes in covariates. These conditions are not sufficient. Since  $h'(z) < 0$ , successive estimates of  $\rho_e$  will oscillate around their limit. In general, a requirement for general convergence of  $\rho_e$  restricts the type of interventions which can be in place. A simple scenario in which  $\rho_e$  cannot converge is provided in Supplementary Section 3.5, and we illustrate an example showing convergence and divergence of  $\rho_e$  in Figure 3. Code to generate a web app that illustrates the problem in general is included in the Supplementary Code.

We may hope that naive updating, when it converges, may solve the optimisation problem in Section 2.3. It does not, and we give a specific counterexample in Supplementary Section 3.4. Finally, we note that the dynamics above also model a related setting, where samples are tracked across epochs and interventions are permanent (Supplementary Section 2). In summary, naive updating can readily lead to wide oscillation of successive risk estimates, and even  $\rho_e$  does converge it is not generally to any useful limit.

## 4 Strategies to avoid this problem

Naive updating is an appropriate method for updating risk scores if no interventions are being made (that is,  $g^a(\rho, x^a) = x^a$  and  $g^\ell(\rho, x^\ell) = x^\ell$ ), as may be the case if a risk score is used for prognosis only, rather than to guide actions<sup>5</sup>. It may also be appropriate if we

<sup>5</sup>EUROscore2 [Nashef et al., 2012] (a risk predictor for cardiac surgery) can be used in this way, by giving patients prognostic estimates but without being used to recommend for or against surgery

do aim to solve the constrained optimisation problem in Section 2.3, and are only concerned with accuracy of the model: in that case, under at least the conditions of Theorem 1, naive updating will lead to estimates  $\rho_e(x^s, x^a)$  converging as  $e \rightarrow \infty$  to a setting in which  $\rho_e$  accurately estimates its own effect: conceptually,  $\rho_e(x^s, x^a)$  estimates the probability of  $Y$  *after* interventions have been made on the basis of  $\rho_e(x^s, x^a)$  itself. Naive updating is otherwise generally not advisable, although a range of alternative modelling strategies do not lead to the same problems.

#### 4.1 More complex modelling and more data

An obvious way to avoid the problem is to model the setting completely, including the effect of any interventions. Methods of this type would include explicit causal modelling, as used in related problems [Sperrin et al., 2018], or counterfactual inference, which has been suggested as a direct approach to the problem [Sperrin et al., 2019]. These approaches would require knowledge or accurate inference of  $g^\ell$  and  $g^a$ , or observation of covariates at several points in each epoch [Sperrin et al., 2018].

A second approach is to consider data from previous epochs alongside the current data when fitting  $\rho_e$ . Such data can be used as a prior on the fitted model [Alaa and van der Schaar, 2018] and could be used to infer model elements:  $\mu_e$ ,  $g^\ell$ ,  $g^a$ , and  $f$ . If accurate data were available, oscillatory effects could even be detected and avoided. A difficulty with this approach in a realistic setting is in distinguishing whether inaccuracies in older models are due to drift in the underlying system [Quionero-Candela et al., 2009] (in our case,  $f$  and  $\mu_e$ ) or due to the effects of intervention. Indeed, the problems with naive updating can be seen as treating model inaccuracies as though they are due to the first effect, when they are in fact due to the second. Definitive assertion of the cause of inaccuracies will, again, generally require more frequent observation of covariates.

#### 4.2 Hold out set

A straightforward and potentially practical means to avoid the problems associated with naive updating is to retain a set of samples in each epoch for which  $\rho_e$  is not calculated, and hence cannot guide intervention. For such samples,  $X_e(0) = X_e(1)$ , so a regression of  $Y$  on  $X_e(0)$  restricted to these ‘held out’ samples can be used as an unbiased estimate for  $f_e$ . If the hold out set is randomly selected, this would emulate a *clinical trial* which enables us to assess the effect of predictive scores (and their associated interventions) across epochs.

A problem with this approach is that any benefit of the

risk score-guided intervention is lost for individuals in the hold-out set. Careful consideration of the ethical consequences of this strategy is therefore required.

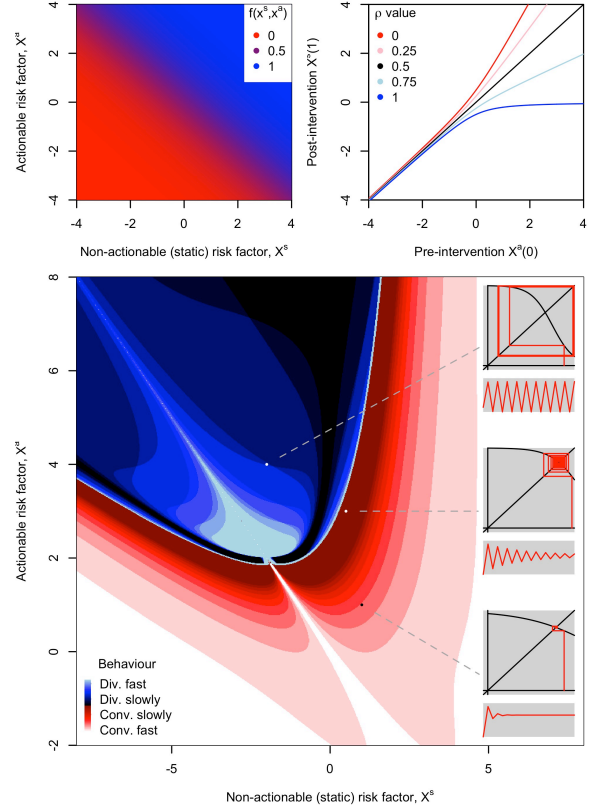


Figure 3: Example showing convergence and divergence of  $\rho_e$  across epochs. We disregard  $x^\ell$ ,  $g^\ell$  in this example. We choose  $f(x^s, x^a) = \text{logit}(x^s, x^a)$  (top left). We choose  $g^a$  with the rationale that we intervene by lowering  $X^a(0)$  when  $\rho_e > 1/2$ , but allow  $X^a(0)$  to increase when  $\rho_e < 1/2$  (that is, resources for intervention are redistributed rather than introduced), and assume that we can intervene more effectively when  $X^a(0)$  is high (strictly,  $g^a(\rho, x^a) = \frac{1}{2} \left( (3 - 2\rho)x^a + (1 - 2\rho)\sqrt{1 + (x^a)^2} \right)$ , top right panel). Bottom panel shows whether  $\rho_e(x^s, x^a)$  converges or diverges, and how long it takes (num. epochs until  $\Delta_e \triangleq |\rho_e - \rho_{e-1}| < 0.01$  or  $(|\Delta_e| > 0.05 \cup |\Delta_e - \Delta_{e-1}| < 0.01)$ ;  $|e| \leq 10$ ). Insets show cobweb plots for relevant recursions, and plots of  $\rho_e$ .

#### 4.3 Control interventions

A radically different option is the direct specification of the interventions  $g_e^\ell$  and  $g_e^a$  in each epoch, considering  $\rho_e$ ,  $\mu_e$  constant, and  $f_e$  to change only slightly with  $e$ . This enables directly addressing the constrained optimisation problem in Section 2.3.

If  $X^\ell$  can be disregarded, and we may regard  $f_{e-1}$  as

an unbiased estimate of  $f_e$ <sup>6</sup>, then we may take a simple inductive approach:

1. At the end of epoch 0, infer  $f_0$  and  $\mu_0$ . Given some fixed functions  $\rho$ ,  $c^a$ , find a function  $g_1^a$  which solves the constrained optimisation problem in section 2.3 assuming  $f_1 = f_0$ ,  $\rho_1 = \rho_0$ . Implement this intervention.
2. At the end of epoch  $e > 0$ , regress  $Y_e$  on

$$X_e(1) = \left( X_e^s(0), g_e^a \left( \rho(X_e^s(0), X_e^a(0)), X_e^a(0) \right) \right)$$

to attain an unbiased estimate of  $f_e$ . Now solve the constrained optimisation problem to optimise  $g_{e+1}^a$ , assuming  $f_{e+1} = f_e$  and  $\rho_{e+1} = \rho_e$

Thus in each epoch an unbiased update of  $f_e$  can be made, and the constrained optimisation problem can be directly solved. If  $X^\ell$  is present, the problem is more complex. We suggest this general case as an open problem (see Supplementary Section 4).

A problem with this approach in a medical setting is that specification of  $g_e^a$  may cause the procedure to be subject to medical device regulation [MHRA, 2019]. Implications of these regulatory processes map to our potential solutions; for example, countries in the EU [EU Council, 2014] have only developed regulatory processes to the point of accommodating static risk scores, and by extension currently treat updated scores as new tools. In these cases a separate evaluation exercise, such as testing on a hold-out, is necessary to demonstrate efficacy prior to dissemination, which would also remedy the problems of naive updating (although costs of repeated formal evaluations of effectiveness, and the ethics of a hold-out, may be a concern). However, the US FDA have proposed an alternative ‘total-life-cycle’ approach [USFDA et al., 2019] which allows for model updating (contingent on defining a performance monitoring mechanism), which, given the problems of naive updating, is potentially seriously flawed.

## 5 Formulation as control-theoretic/reinforcement learning problem

Control theory [Bertsekas, 1995] and its modern incarnation, reinforcement learning [Sutton and Barto, 2018], study temporal problems where multiple actions are available at each time step. The aim of the field is to come up with an optimal policy either from the start or, in the partially observable case, a mechanism that

<sup>6</sup>This assumption underlies the fundamental point of a risk score

quickly converges to the optimal policy. In the latter the regret is considered to be how much utility is lost compared to using the optimal policy from the start. The methods underlying this, like dynamic programming, are used in a variety of fields such as; playing go [Silver et al., 2018], in dynamic treatment strategy [Alaa and van der Schaar, 2018] and mechanical and electrical engineering. Here we use the formulation of a Partially Observable Markov Decision Processes (POMDP) [Yuksel, 2017], and adopt the notation from [Wang et al.] whereby we consider the POMDP as a 7-Tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \Omega, \mathcal{Z}, \gamma)$ :

- $\mathcal{S}, \mathcal{A}$  and  $\Omega$  are spaces of states, actions and observations.
- $\mathcal{T}$  is the transition kernel that describes the evolution given state and action, e.g.  $s_{e+1} \sim \mathcal{T}(\cdot | s_e, a_e)$  (i.e. a set of conditional transition probabilities between states and actions).
- $\mathcal{Z}$  is a kernel for the observation given the state, e.g.  $o_{e+1} \sim \mathcal{Z}(\cdot | s_e, a_e)$ <sup>7</sup>.
- $r_e$  represents our reward for being in state  $s$  and taking action  $a$  at time (or equivalently epoch)  $e$ , and is sampled from  $\mathcal{R}$  - i.e.  $r_e \sim \mathcal{R}(s_e, a_e)$
- $\gamma$  is a discount factor that down-weights future rewards if  $0 < \gamma < 1$ .

A solution candidate is a policy

$$a_e \sim \pi \left( \{o_s, r_s, a_s\}_{s=1}^{e-1} \right)$$

which aims to maximise

$$\mathbb{E} \sum_{e=1}^M \gamma^{e-1} r(s_e, a_e)$$

where  $M$  represents the maximum number of time/epoch steps. Other reward/utility parametrisations are possible e.g. to include a final pay off or infinite time horizon pay off. The beauty of this framework is the flexibility: aspects such as optimisation under uncertainty can be included by including parameters of reward, transition and observation processes into the (unobserved) state variable.

We cast the above in this framework:

$$\begin{aligned} s_e &= (X_e(0), X_e(1), Y_e) \\ a_e &= \rho_e \\ o_e &= ((X_e^s(0), X_e^a(0)), Y_e) \\ r_e &= \mathbb{P}(\bar{Y}_{e+1} | s_e, a_e) \end{aligned}$$

<sup>7</sup>Note that here future observations depend on current states and actions and not on future states and actions



The transition kernel from  $s_e$  to  $s_{e+1}$  consists of; sampling  $X_{e+1}(0)$  (note that this sampling is independent of  $s_e$ ), intervening using this sample with  $\rho_e$  to form  $X_{e+1}(1)$ , and then using these values to sample  $Y_{e+1}$  from the resulting conditional distribution. Finally we note that given Assumption 5 our policy  $a_e \sim \pi(o_e, r_e, a_e)$  as previous epochs are ignored. Indeed, this assumption also implies that  $s_{e+1}$ ,  $o_{e+1}$  and  $r_e$  only depend on the previous state through  $a_e = \rho_e$ . In the control view point it is also easy to formulate the longitudinal problem (this corresponds to setting  $X_{e+1}(0) = X_e(1)$ ).

The description above allows to use methods of the field such as Q-learning, (approximate dynamic programming), PDE-based approaches such as the Hamilton Jacobi Bellman equation and many more. These methods create a policy which maps the historical observations to an action (for the problem at hand a risk score function). Most of the rigorous methods require a low dimensional state space [Powell, 2007].

## 6 Discussion

In this work, we elaborate on the issue raised by Lenert and Sperrin [Lenert et al., 2019, Sperrin et al., 2019] and propose a framework for quantitatively modelling its effects, with a particular focus on a model which is updated repeatedly. We demonstrate some consequences of ignoring this problem, and note that they occur even in highly idealised circumstances. Although the problem can generally be avoided by more complex and complete modelling, we consider that this is often impractical: a full consideration of the setting in which a model will eventually be used is not generally considered until the model is to be implemented [Lipton and Steinhardt, 2018].

The formulation of the constrained optimisation problem in section 2.3 makes it clear that for fixed  $g^\ell$ ,  $g^a$ , the best possible  $\rho_e$  is not necessarily the oracle estimator in equation 3. However, many machine learning models tend to focus on accurate prediction of outcomes [Nashef et al., 2012], rather than directly solving problems of the type in section 2.3; hence, the naive updating setting considers a  $\rho_e$  which does exactly this. In the naive updating setting, we are assuming an analyst who ignores this effect.

The model presented here is not a full description of modern predictive scoring systems; however, it is extensible in various ways (some detailed in Supplementary Section 4). In particular,  $g^\ell$  and  $g^a$  could be random-valued rather than deterministic. We also note that we assume a covariate value after intervention confers the same contribution to risk of  $Y$  as it does when it takes the same value ‘naturally’, which

may not be realistic.

We assume we are ‘starting over’ with new samples at the beginning of each epoch, and for naive updating, we assume that covariate values are identically distributed. The basis for this assumption is that we generally expect interventions to be zero-sum: that is, the risk score guides a redistribution of intervention rather than introduction of interventions, so the total effect on the sample population remains roughly the same in each epoch. In this assumption, we differ from that in the analysis by Lenert [2019]. We can alternatively interpret this assumption as taking all interventions as being short-term and having ‘worn off’ by the start of the next epoch. The problem raised here also exists for the more general setting when interventions have long term effects and we consider longitudinal effects.

In the setting where models change at each epoch, if  $m_{\tilde{f}_e}$  is known at the current epoch  $e$ , we note a fair comparison of models is one which compares models built using the training data available at the current epoch<sup>8</sup>. If  $m_{\tilde{f}_e}$  is not known, then a holdout set for test data must be used so a fair comparison can be made using an estimate of  $m_{\tilde{f}_0}$  (assuming  $\tilde{f}_0 \approx f$ ). This is because at epoch  $e$  we only have access to  $(X_e(0), Y_e)$  and not  $X_e(1)$ , and so we are not able to properly gain insight to the behaviour of  $\tilde{f}_e$  needed to provide an estimate of  $m_{\tilde{f}_e}$ . An attempt to estimate  $m_{\tilde{f}_e}$  using  $(X_e(0), Y_e)$  implicitly assumes that  $Y_e$  directly depends on  $X_e(0)$ , and as a result  $\rho_e$  would appear much closer to  $\tilde{f}_e$  than is the case. Put simply, by implementing naive model updating not only may performance severely worsen (even if better models were used), but in not providing a holdout test set stakeholders may not even be able to recognise that performance is worsening as the number of epochs increase.

In essence, we provide a causal framework within which to understand a crucial issue in regulation of machine learning and AI-based tools in health and further afield, demonstrating that approaches which incorporate naive updating are unlikely to be fit for purpose. Moreover, even where solutions are available to address the bias introduced by updating on ‘real-world’ data in which outcomes represent (at least in part) the effects of an algorithm, these restrict the potential of ‘online’ and frequently updated solutions. We hope that our work will foster discussion of this interesting problem, which is becoming increasingly pertinent as machine-learning based predictive scores become widely used to guide decision making, and policymakers act to address how to regulate these tools to ensure safety and effectiveness.

<sup>8</sup>This is not to say that the performance of models will not deteriorate over epochs, just that the issue may not lie with the model structure.

## References

- A. M. Alaa and M. van der Schaar. Autoprognosis: Automated clinical prognostic modeling via bayesian optimization with structured kernel learning. *arXiv preprint arXiv:1802.07207*, 2018.
- D. P. Bertsekas. *Dynamic programming and optimal control*, volume 1. Athena scientific Belmont, MA, 1995.
- EU Council. EU regulation no 2017/745 on medical devices, 2014. <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32017R0745>.
- J. Friedman, T. Hastie, and R. Tibshirani. *The Elements of Statistical Learning*, volume 1. Springer Series in Statistics New York, 2001.
- S. L. Hyland, M. Faltys, M. Hüser, X. Lyu, T. Gumbusch, C. Esteban, C. Bock, M. Horn, M. Moor, B. Rieck, et al. Early prediction of circulatory failure in the intensive care unit using machine learning. *Nature Medicine*, 26(3):364–373, 2020.
- M. C. Lenert, M. E. Matheny, and C. G. Walsh. Prognostic models will be victims of their own success, unless.... *Journal of the American Medical Informatics Association*, 26(12):1645–1650, 2019.
- Z. C. Lipton and J. Steinhardt. Troubling trends in machine learning scholarship. *arXiv preprint arXiv:1807.03341*, 2018.
- MHRA. Medical device stand-alone software including apps (including IVDMDs), 2019.
- S. A. Nashef, F. Roques, L. D. Sharples, J. Nilsson, C. Smith, A. R. Goldstone, and U. Lockowandt. Euroscore ii. *European Journal of Cardio-Thoracic Surgery*, 41(4):734–745, 2012.
- W. B. Powell. *Approximate Dynamic Programming: Solving the Curses of Dimensionality*. John Wiley & Sons, Oct. 2007.
- J. Quionero-Candela, M. Sugiyama, A. Schwaighofer, and N. D. Lawrence. *Dataset Shift in Machine Learning*. The MIT Press, 2009.
- F. Rahimian, G. Salimi-Khorshidi, A. H. Payberah, J. Tran, R. A. Solares, F. Raimondi, M. Nazarzadeh, D. Canoy, and K. Rahimi. Predicting the risk of emergency admission with machine learning: Development and validation using linked electronic health records. *PLoS Medicine*, 15(11):e1002695, 2018.
- Z. R. Shi, Z. S. Wu, R. Ghani, and F. Fang. Bandit data-driven optimization: AI for social good and beyond. Aug. 2020.
- D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362(6419):1140–1144, 2018.
- M. Sperrin, G. P. Martin, A. Pate, T. Van Staa, N. Peek, and I. Buchan. Using marginal structural models to adjust for treatment drop-in when developing clinical prediction models. *Statistics in Medicine*, 37(28):4142–4154, 2018.
- M. Sperrin, D. Jenkins, G. P. Martin, and N. Peek. Explicit causal reasoning is needed to prevent prognostic models being victims of their own success. *Journal of the American Medical Informatics Association*, 26(12):1675–1676, 2019.
- R. S. Sutton and A. G. Barto. *Reinforcement Learning, second edition: An Introduction*. MIT Press, Nov. 2018.
- USFDA et al. Proposed regulatory framework for modifications to artificial intelligence/machine learning (AI/ML)-based software as a medical device (samd)-discussion paper, 2019. <https://www.fda.gov/files/medical%20devices/published/US-FDA-Artificial-Intelligence-and-Machine-Learning-Discussion-Paper.pdf>.
- E. Wallace, E. Stuart, N. Vaughan, K. Bennett, T. Fahy, and S. M. Smith. Risk prediction models to predict emergency hospital admission in community-dwelling adults: a systematic review. *Medical Care*, 52(8):751, 2014.
- Y. Wang, B. Liu, J. Wu, Y. Zhu, S. S. Du, L. Fei-Fei, and J. B. Tenenbaum. DualSMC: Tunneling differentiable filtering and planning under continuous POMDPs. *ijcai.org*.
- S. Yuksel. Control of stochastic systems. *Queen's University Mathematics and Engineering and Mathematics and Statistics*, 2017.

---

# Model updating after interventions paradoxically introduces bias

## Supplementary Materials

---

### 1 Example of functions and variables in a realistic setting

We consider the model proposed by Rahimian et al Rahimian et al. (2018) for prediction of emergency admission to a hospital in a given time period on the basis of electronic health records (EHRs). Such a model is not in common use in the location considered (England), so the data in the original paper is not affected by the problems we describe in the main manuscript.

For clarity<sup>1</sup>, we presume a prediction window of ten months (February-November), and that predictions are made and distributed to primary health practitioners in January, with a new model being trained on the basis of each year's data in December, to be implemented the following January. In this setting, distribution of the score may open a second causal pathway between covariates and outcome as shown in figure 1, and is thus susceptible to the problems of naive updating.

In this setting, variables and functions may be interpreted as follows:

1.  $Y$  the event 'an emergency admission in the following year'
2.  $X_e(0)$  the values of all variables which affect  $E(Y)$  at the time when the predictive score is computed (the start of each year)
3. An 'epoch': the time in which a given model is in use; eg, each year.
4. 'Time':  $t = 0$  when the predictive score is computed (the start of January);  $t = 1$  represents the time after which any interventions are made (the start of February).
5.  $X_e^S$  covariates affecting  $\mathbb{E}(Y)$  which are included in the predictive score but which cannot be directly modified in the time frame: age, time since most recent emergency admission
6.  $X_e^A$  covariates affecting  $\mathbb{E}(Y)$  included in the predictive score which can be modified in the time frame: current medications.
7.  $X_e^L$  covariates affecting  $\mathbb{E}(Y)$  which are not included in the predictive score, and possibly can be modified in the time frame: blood pressures, cardiac function
8.  $f_e$  the underlying causal process for  $Y$  given patient status; that is, the probability of admission in the subsequent year, given
9.  $g_e^A$  Hypothetical prescribed interventions made on  $X^A$  in response to a predictive score; for instance, reduce drug dosages. We roughly assume that this intervention is symmetric; for a patient at low emergency risk, a higher drug dose is acceptable.
10.  $g_e^L$  Hypothetical prescribed interventions made on  $X^L$  in response to a predictive score; for instance, treat low or high blood pressure.

It is clear that if such a risk score were used universally, and data was collected from the period in which a model was in place was then, then the data would be affected by the effect of the predictive score itself.

The model does not fully describe this setting. The trichotomisation into  $X^L$ ,  $X^A$ , and  $X^S$  is not perfect; intervention on  $X^L$  could also affect some variables in  $X^A$  and vice versa. Interventions are likely to be random-valued to some extent.

---

<sup>1</sup>Analogous times and variables can be described for other prediction periods and updating patterns

## 2 Alternative system described by naive updating

We note that the definition of  $h$  (equation 9), and hence the following comments on recursion dynamics, can be used to describe a related setting in which we track the same samples over epochs, and the effect of interventions  $g^a, g^\ell$  remain in place. Formally, we retain definitions of  $X^s, X^a, X^\ell, e, t, f_e, g_e^a, g_e^\ell, \rho_e$  and all assumptions except 4,7. In place, we assume that  $f_e, g_e^a, g_e^\ell$  are fixed across epochs, but instead of resampling  $X_e(0)$  from  $\mu_e$ , we have

$$X_{e+1}(0) = X_e(1) \quad (1)$$

thus, while values  $X_0(0)$  are sampled from the distribution  $\mu_0$ , values  $X_e(0)$  are then determined for  $e > 0$ . We illustrate this in figure 1

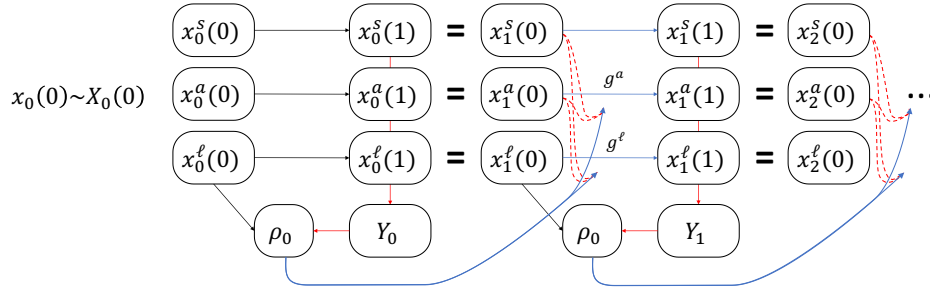


Figure 1: Diagram showing alternative setup for naive updating. Values  $x^s, x^a, x^\ell$  are sampled at  $(e, t) = (0, 0)$ , and used to determine  $\rho_0$ . Values are conserved until  $t = 1$ , and remain the same at the start of epoch 1  $((e, t) = (1, 0))$ . Values are intervened on by  $g^a, g^\ell$  according to  $\rho_0(x_1^s(0), x_1^a(0))$ , and resultant values at  $(e, t) = (1, 1)$  are conserved until the start of the next epoch at  $(e, t) = (2, 0)$ . Lowercase letters indicates that, while quantities random-valued, they inherit all randomness from their values at  $(e, t) = (0, 0)$ . Colour and line conventions are as for figure 2

Now formulas 8, 9 will hold, and the recursion will proceed as detailed in theorem 1.

## 3 Proofs and counterexamples

### 3.1 Optimising both $\rho$ and $g^a, g^\ell$ is equivalent to a general resource allocation problem

Considering the constrained optimisation problem in section 2.3. We show that if we allow  $\rho$  and  $g^a, g^\ell$  to vary independently, then the constrained optimisation is equivalent to the solution of a problem in which the use of a predictive score is redundant.

**Theorem 1.** Suppose that the triple  $(\rho_{opt}, g_{opt}^a, g_{opt}^\ell)$  minimises quantity 4 subject to constraint 5 in section 2.3, where all are arbitrary functions of two variables in the appropriate range. Let  $h_{opt}^a$  and  $h_{opt}^\ell$  be solutions to a second constrained optimisation problem: find  $h^a(x^s, x^a)$  and  $h^\ell(x^s, x^a, x^\ell)$  which minimise

$$\begin{aligned} \mathbb{E}_{X_e(0)} \{ & f(X^s, \\ & h^a(X_e^s(0), X_e^a(0)), \\ & h^\ell(X_e^s(0), X_e^a(0), X_e^\ell(0))) \} \end{aligned} \quad (2)$$

subject to

$$\begin{aligned} \mathbb{E}_{X_e(0)} \{ & c^a(X_e^a(0), \\ & X_e^a(0) - h^a(X_e^s(0), X_e^a(0))) + \\ & c^\ell(X_e^\ell(0), \\ & X_e^\ell(0) - h^\ell(X_e^s(0), X_e^a(0), X_e^\ell(0))) \} \leq C \end{aligned} \quad (3)$$

with  $c^a, c^\ell, f$  as for section 2.3.

Then the minima of quantity 4 in the main text and of quantity 2 achieved by  $(\rho_{opt}, g_{opt}^a, g_{opt}^l)$  and  $(h_{opt}^a, h_{opt}^l)$  are the same.

*Proof.* Given a triple  $(\rho_{opt}, g_{opt}^a, g_{opt}^l)$ , we explicitly construct an  $(h_{opt}^a, h_{opt}^l)$  which attains the same minimum, and vice versa.

Given  $(\rho_{opt}, g_{opt}^a, g_{opt}^l)$ , the corresponding forms of  $h_{opt}^a, h_{opt}^l$  are simply

$$\begin{aligned} h_{opt}^a(x^s, x^a) &= g_{opt}^a(\rho(x^s, x^a), x^a) \\ h_{opt}^l(x^s, x^a, x^l) &= g_{opt}^l(\rho(x^s, x^a), x^l) \end{aligned} \quad (4)$$

Given  $h_{opt}^a, h_{opt}^l$ , the correspondence is slightly more complex. Set  $\rho_{opt}$  as a bijective function from  $\mathbb{R}^{n_s+n_a}$  to  $\mathbb{R}$ ; for instance, set it to ‘splice’ the decimal digits of arguments together. Now set  $g_{opt}^a, g_{opt}^l$  to firstly ‘decrypt’ the value of  $\rho_{opt}$  back into constituent parts ( $x^s$  and  $x^a$ ), and then compute  $h_{opt}^a(x^s, x^a)$  and  $h_{opt}^l(x^s, x^a, x^l)$  as outputs.

This shows that the two constrained optimisation problems are equivalent.  $\square$

We note that this implies that optimising  $(\rho, g^a, g^l)$  jointly is equivalent to a more general treatment-allocation problem which does not involve a predictive score.

### 3.2 Counterexample showing naive updating can cause better models to appear worse

For this counterexample we shall use the following set up:

$$f(x^s, x^a, x^l) = f(x^s, x^a) = (1 + e^{-x^s - x^a})^{-1} \quad (5)$$

$$\rho_0(x^s, x^a \mid X_0^*, Y_0^*) = \begin{cases} \frac{\sum_{i=1}^n (Y_0^*)_i \mathbb{1}\{\sum_{j=1}^2 (X_0^*)_{ij} > 0\}}{\sum_{i=1}^n \mathbb{1}\{\sum_{j=1}^2 (X_0^*)_{ij} > 0\}} & x^s + x^a > 0 \\ \frac{\sum_{i=1}^n (Y_0^*)_i \mathbb{1}\{\sum_{j=1}^2 (X_0^*)_{ij} \leq 0\}}{\sum_{i=1}^n \mathbb{1}\{\sum_{j=1}^2 (X_0^*)_{ij} \leq 0\}} & x^s + x^a \leq 0 \end{cases} \quad (6)$$

$$\rho_1(x^s, x^a \mid X_1^*, Y_1^*) = (1 + e^{-\hat{\beta}_0 - x^s \hat{\beta}_1 - x^a \hat{\beta}_2})^{-1} \text{ where } \hat{\beta} = \operatorname{argmax}\{\mathcal{L}(\beta \mid X_1^*, Y_1^*)\} \quad (7)$$

$$m_{\tilde{f}_e}(\rho_e \mid X_e^*, Y_e^*) = \mathbb{E}_\mu [|f(X^s, g^a(\rho_{e-1}, X^a)) - \rho_e(X^s, X^a \mid X_e^*, Y_e^*)|] \quad (8)$$

$$g^a(\rho, x^a) = (1 - \rho)(x^a + 3) + \rho(x^a - 3) \quad (9)$$

For simplicity, we shall view the latent variables as having no effect on the true risk score  $f$ , which corresponds to the scenario where (if no interventions are made), it is possible with the data we observe to fully specify  $f$ . For the purpose of the counterexample it is reasonable to do this as model performance only requires  $m_{\tilde{f}_e}$ , which has no dependence on latent covariates.

We also state, that due to the omission of latent covariates,  $X_e(0) = (X_e^s(0), X_e^a(0)) \sim N_2(0, I_2)$ , which is then used to generate (through the statistical program R) an initial training data set at epoch 0, of size  $n = 100$ , which is summarised below:

index	$(\mathbf{X}_0^*)_{.1}$	$(\mathbf{X}_0^*)_{.2}$	$\mathbf{Y}_0^*$
1	1.185	1.272	1
2	0.881	-0.995	0
3	0.122	-0.956	0
$\vdots$			
98	-0.826	1.779	1
99	0.853	0.151	1
100	0.177	0.805	1

This training data can then inputted into  $\rho_0$  to give the following function:

$$\rho_0(x^s, x^a \mid X_0^*, Y_0^*) = \begin{cases} 0.733 & x^s + x^a > 0 \\ 0.200 & x^s + x^a \leq 0 \end{cases} \quad (10)$$

When intervening on any covariates at epoch 1 the function given in equation 10 will be used to produce  $X_1(1)$  and subsequently  $Y_1$ .

We now consider  $\mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)]$ , which we approximate using a Monte Carlo estimate with 1000 samples. However,  $m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)$  also requires approximation, and so a Monte Carlo estimate with the same number of samples is also used for this function. The procedure is as follows:

1. For  $i$  from 1 to 1000:
  - (a) Obtain a training data set,  $(X_0^*, Y_0^*)_i$ , by taking  $n$  samples of  $(X_0(0), Y_0)$ .
  - (b) Use this training data set to obtain a  $(\rho_0)_i$  of the form given in equation 10.
  - (c) For  $j$  from 1 to 1000:
    - i. Sample  $(x^s, x^a)_j \sim X_0(0)$ .
  - (d)  $m_{\tilde{f}_0}(\rho_0 | (X_0^*, Y_0^*)_i) \approx \frac{1}{1000} \sum_{j=1}^{1000} |f((x^s, x^a)_j) - \rho_0((x^s, x^a)_j | (X_0^*, Y_0^*)_i)|$
2.  $\mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)] \approx \frac{1}{1000} \sum_{j=1}^{1000} m_{\tilde{f}_0}(\rho_0 | (X_0^*, Y_0^*)_i)$

With this in mind, we give the following approximation:  $\mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)] \approx 0.124$ .

If we assert that interventions never take place, then we can use the same procedure described above to obtain  $\mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_1 | X_0^*, Y_0^*)] \approx 0.056$ . So here we can clearly see that in the setting where interventions are never made,  $\mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)] > \mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_1 | X_0^*, Y_0^*)]$ , and so the model closer to the truth is the logistic regression model at epoch 1. If agents were allowed to make interventions (based on (10)) however, we would consider  $\mathbb{E}_{(X_1^*, Y_1^*)} [m_{\tilde{f}_1}(\rho_1 | X_1^*, Y_1^*)] \approx 0.197$  instead. Now, since  $\mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)] < \mathbb{E}_{(X_1^*, Y_1^*)} [m_{\tilde{f}_1}(\rho_1 | X_1^*, Y_1^*)]$ , we would come to the incorrect conclusion that the model closer to the truth is the model used at epoch 1. Consequently we can state that, given the setup provided in section 3.1,

$$\begin{aligned} \mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)] &> \mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_1 | X_0^*, Y_0^*)] \not\Rightarrow \\ \mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)] &> \mathbb{E}_{(X_1^*, Y_1^*)} [m_{\tilde{f}_1}(\rho_1 | X_1^*, Y_1^*)] \end{aligned} \quad (11)$$

Additionally, we show that for this example:

$$\begin{aligned} \mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)] &> \mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_1 | X_0^*, Y_0^*)] \not\Rightarrow \\ \mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)] &> \mathbb{E}_{(X_1^*, Y_1^*)} [m_{\tilde{f}_0}(\rho_1 | X_1^*, Y_1^*)] \end{aligned} \quad (12)$$

as  $\mathbb{E}_{(X_1^*, Y_1^*)} [m_{\tilde{f}_0}(\rho_1 | X_1^*, Y_1^*)] \approx 0.215 > 0.124 \approx \mathbb{E}_{(X_0^*, Y_0^*)} [m_{\tilde{f}_0}(\rho_0 | X_0^*, Y_0^*)]$ . This statement is given here because for  $\tilde{f}_0$ , and therefore  $m_{\tilde{f}_0}$ , it is possible to gain estimates through a holdout test data set. Whilst the comparison is not between a risk score ( $\rho_e$ ) and the function it is trying to estimate ( $\tilde{f}_e$ ), the effect of deteriorating performance as epochs increase is still captured. Going further, it is assumed that if stakeholders were implementing naive model updating, they would assume that  $\rho_e$  is estimating  $\tilde{f}_0$  for all epochs as the belief is that interventions do not effect the model. Therefore, comparison with  $\tilde{f}_0$  will heighten the impression to stakeholders that using an updated model structure is causing performance to deteriorate, especially for epoch 0 to epoch 1, where for this comparison  $\rho_0$  is actually estimating  $\tilde{f}_0$ .

We expect from a stakeholders view that comparison (using estimates) between the two models at successive epochs usually leads to the inequality  $m_{\tilde{f}_0}(\rho_{e-1} | X_{e-1}^*, Y_{e-1}^*) < m_{\tilde{f}_0}(\rho_e | X_e^*, Y_e^*)$ , and therefore the conclusion is that the new model leads to worse performance. We advise that a conclusion is only reached after further comparison is done between  $m_{\tilde{f}_0}(\rho_{e-1} | X_e^*, Y_e^*)$  and  $m_{\tilde{f}_0}(\rho_e | X_e^*, Y_e^*)$ , as this gives an indication whether the drop in performance is due to the model structure or the intervention effect.

Finally, we advise caution when considering the effect of latent variables when estimating  $m_{\tilde{f}_0}(\rho_e | X_e^*, Y_e^*)$ . This is due to that fact that when holdout test data is used to obtain an estimate, it is an estimate of  $f$  rather than

an estimate of  $\tilde{f}_0$ . If the latent variables have a small influence on  $f$  than  $f \approx \tilde{f}_0$  and we can make inferences as shown above, but if latent variables have a large influence on  $f$  then our comparison is not based on  $m_{\tilde{f}_0}$  but instead on  $m_f$ . This creates a problem as now how well we perceive our model's performance can be determined largely by how well a model arbitrarily captures the latent covariate information using just the set and actionable covariates. It therefore becomes substantially more difficult to determine whether the cause of a model's poor performance is due to the model, the intervention effect or insufficient data. As a general rule however, large values of  $m_{\tilde{f}_0}(\rho_0|X_0^*, Y_0^*)$  should indicate that either the initial model is very poor or that there is insufficient data, but in either case careful consideration of what could possibly influence the underlying mechanism should be made before a risk score is built and given to agents, to ensure that latent variables affect the model as little as possible.

### 3.3 Proof of theorem 1

If  $h'(z_0) \leq -1$  then the single fixed point of  $h$  is unstable and  $\rho_n$  cannot converge to it unless it was always equal to  $z_0$ . There can be no other  $z$  with  $h(z) = z_0$  since  $h'(z) < 0$  by assumption.

Since  $\rho_e \in [0, 1]$  and  $h'(z) < 0$ ,  $\rho_e$  must converge to either an oscillation between two values, or to a single value.

If the bounds on partial derivatives hold, then from the triangle and Cauchy-Schwarz inequalities, for  $z \in R$

$$\begin{aligned}
|h'(z)| &\leq \mathbb{E}_{X^L} \left[ \sum_i^{p^A} |\delta_i^{g^A} \delta_i^{f^A}| + \sum_i^{p^L} |\delta_i^{g^L} \delta_i^{f^L}| \right] \\
&= \sum_i^{p^A} |\delta_i^{g^A}| \mathbb{E}_{X^L} [|\delta_i^{f^A}|] + \sum_i^{p^L} \mathbb{E}_{X^L} [|\delta_i^{g^L} \delta_i^{f^L}|] \\
&\leq \sqrt{\sum_i^{p^A} (\delta_i^{g^A})^2 \sum_i^{p^A} \mathbb{E}_{X^L} [\delta_i^{f^A}]^2} \\
&\quad + \sqrt{\sum_i^{p^L} \mathbb{E}_{X^L} [(\delta_i^{g^L})^2] \sum_i^{p^L} \mathbb{E}_{X^L} [(\delta_i^{f^L})^2]} \\
&\leq \sqrt{k_1 k_3} + \sqrt{k_2 k_4} < 1
\end{aligned} \tag{13}$$

so the map  $h : \rho_n \rightarrow \rho_{n+1}$  is a contraction, and the convergence of the recurrence  $\rho_n \rightarrow \rho_{n+1}$  follows from the Banach fixed-point theorem, as long as  $\rho_n \in R$  for some value of  $n$ .

### 3.4 Counterexample showing failure of naive updating to generally solve constrained optimisation problem

For this counterexample, we do not need to consider latent covariates, and will assume they do not exist.

Under the setting in section 2.2, if  $\rho_n$  converges to  $\rho_\infty(x^s, x^a)$  for some  $x^s, x^a$  under naive updating, then we have

$$\rho_\infty(x^s, x^a) = h(\rho_\infty(x^s, x^a)) = f(g(\rho_\infty(x^s, x^a), x^a), x^s) \tag{14}$$

Suppose  $x^s$  and  $x^a$  each have dimension 1, and consider the example:

$$\begin{aligned}
f(x^a, x^s) &= \text{logit}(x^a + x^s) = \frac{1}{1 + \exp(-(x^a + x^s))} \\
g(\rho, x^a) &= x^a - \log(1 + \rho) \\
c_A(x) &= x
\end{aligned}$$

For a given function  $\rho$ , the objective and cost are, respectively

$$\begin{aligned}
\text{obj}\{\rho\} &= E \{ (1 + (1 + \rho) \exp(-(X^s + X^a)))^{-1} \} \\
\text{cost}\{\rho\} &= E \{ \log(1 + \rho) \}
\end{aligned} \tag{15}$$

Using an oracle predictor of  $Y|X$ , as in the previous section,  $\rho_n$  converges to the fixed point of the recursion  $z \rightarrow f(g(z, x^a), x^s)$ , which is

$$\rho_\infty(x^s, x^a) = \frac{1}{2} \left( \sqrt{(e^{x+y} + 1)^2 + 4e^{x+y}} - (e^{x+y} + 1) \right) \quad (16)$$

To see why this is not optimal, suppose  $X^a, X^s$  have a discrete distribution taking either of the values  $(0, -1)$ ,  $(0, 1)$  with probability  $1/2$ . Then

$$\begin{aligned} \text{cost}\{\rho_\infty\} &= \frac{\log(2)}{2} \approx 0.346 \\ \text{obj}\{\rho_\infty\} &= \frac{1+e}{1+e+\sqrt{1+6e+e^2}} \approx 0.428 \end{aligned}$$

However, consider some  $\rho_0$  with  $\rho_0(0, -1) = 0$ ,  $\rho_0(0, 1) = 1$ . Now

$$\begin{aligned} \text{cost}\{\rho_0\} &= \frac{\log(2)}{2} = \text{cost}\{\rho_\infty\} \\ \text{obj}\{\rho_0\} &= \frac{1}{2} \left( \frac{1}{1+e} + \frac{e}{2+e} \right) \approx 0.423 < \text{obj}\{\rho_\infty\} \end{aligned} \quad (17)$$

### 3.5 Simple example of updating leading to oscillation

Define  $g(\rho, x^a)$  as above, and instead define

$$f(x^a, x^s) = \text{logit}(-k(x^a + x^s)) \quad (18)$$

As usual, we presume that to estimate  $\rho$ , we regress  $Y$  on  $X_0^S, X_0^A$ , and we do it accurately enough to presume  $\rho$  is an oracle. Now

$$\begin{aligned} h(x) &= \frac{1}{1 + (1+x)^k \exp(-k(x^s + x^a))} \\ h'(x) &= -k \frac{e^{k(x^s + x^a)}(1+x)^{k-1}}{(e^{k(x^s + x^a)} + (1+x)^k)^2} \end{aligned} \quad (19)$$

Consider a setting when  $x^s = x^a = 0$  and  $k = 8$ . Now  $h(0) = 1/2 > 0$  and  $h(1/5) \approx 0.189 < 1/5$ . For  $x \in (0, 1)$  we have  $h'(x) < 0$ , so the equation  $h(x) = x$  has a single solution in  $(0, 1/5)$ . But on  $(0, 1/5)$ , we have  $h'(x) < -1$ . So if  $x_0$  is the unique root of  $h(x) - x$  on  $x \in (0, 1)$  then  $h'(x_0) < 0$

Now as long as  $\rho_0(x^s, x^a)$  is not exactly the value of  $x$  for which  $h(x) = x$ , if we update  $\rho_n$  using  $h$ , it can never converge as the fixed point of the map  $h$  is unstable.

Conceptually, although no intervention changes  $x^a$  very much, the function  $f$  is very sensitive to small changes in  $x^a$  when  $k = 8$ , so a small change in  $x^a$  will necessarily cause a larger change in  $f(x^a, x^s)$  when  $\rho$  is near the fixed point of  $h$ .

## 4 Open problems

We propose the following short list of open problems in this area.

1. Determine a framework to modulate both  $g^L$  and  $g^A$  with the aim of solving the constrained optimisation problem in section 2.3.
2. Determine the dynamics and consequences of other model-updating strategies. What happens if training data is aggregated at each step, rather than only the most recent data being used?
3. Derive results of successive adjuvancy in more general circumstances.
4. How do the dynamics of the model change when assumptions differ? Can  $f$ ,  $g^L$  and  $g^A$  be extended to be random-valued, and possibly agglomerated into a single intervention function?
5. How can assumptions be changed to approximate more general machine learning settings?



---

## References

Rahimian, F., Salimi-Khorshidi, G., Payberah, A. H., Tran, J., Solares, R. A., Raimondi, F., Nazarzadeh, M., Canoy, D., and Rahimi, K. (2018). Predicting the risk of emergency admission with machine learning: Development and validation using linked electronic health records. *PLoS Medicine*, 15(11):e1002695.